

Faculty of Engineering and Information Technology

University of Technology, Sydney

# SOCIAL SECURITY DATA MINING: AN AUSTRALIAN CASE STUDY

A Thesis Submitted in Fulfilment of the Requirements for The Degree of  
Doctor of Philosophy, Faculty of Engineering and Information Technology,  
University of Technology Sydney.

By

Hans Michael Bohlscheid

October 2013

## **CERTIFICATE OF AUTHORSHIP/ORIGINALITY**

I, Hans Michael Bohlscheid, certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of Requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

Signature of Candidate

---

Hans M Bohlscheid

01 /10 /2013

## ACKNOWLEDGEMENT

First and foremost, I would like to express my deepest gratitude to Professor Longbing Cao for his supervision and ongoing support over many years. I regard Professor Cao as an energetic, inspiring individual and will always remember our many discussions and his spontaneous ideas.

I also remain grateful to my colleagues and friends Drs Huaifeng Zhang and Yanchang Zhao for their advice regarding the technical aspects of my thesis.

In the mix of the many people to whom I am indebted are Yuri Zubrutsky, Daniel Marlay and Richard Brookes who independently evaluated my results and last but not least, I extend my heartfelt thanks to my dear colleagues Peter Newbigin and Brett Clark for their expertise in data extraction and their acute knowledge of Centrelink business practices.

In their own unique way, each of the above contributed to the enjoyment and success of my research and collectively, they are responsible for the skills and subject matter knowledge I possess today.

*“I was striking an uneasy balance between the ambition I had for myself, and what those closest to me expected of me. So I stopped pretending to myself that I was anything other than what I was, and began to direct all my energy into finishing the only work that mattered to me.”*

*J.K. Rowling*

# TABLE OF CONTENTS

<b>SOCIAL SECURITY DATA MINING: AN AUSTRALIAN CASE STUDY .....</b>	<b>1</b>
<b>CERTIFICATE OF AUTHORSHIP/ORIGINALITY .....</b>	<b>2</b>
<b>ACKNOWLEDGEMENT .....</b>	<b>3</b>
<b>TABLE OF CONTENTS .....</b>	<b>4</b>
<b>LIST OF FIGURES.....</b>	<b>9</b>
<b>LIST OF TABLES .....</b>	<b>10</b>
<b>ABSTRACT.....</b>	<b>12</b>
<b>1 INTRODUCTION.....</b>	<b>14</b>
1.1 BACKGROUND.....	14
1.2 BUSINESS NEEDS AND RESEARCH ISSUES .....	15
1.2.1 <i>Business Needs : Social Security Service Delivery, Reform and Challenge .....</i>	<i>15</i>
1.2.2 <i>Research Issues in Social Security Data Mining .....</i>	<i>17</i>
1.3 RESEARCH METHODOLOGY .....	19
1.3.1 <i>Positive / Negative Sequential Rules Mining for Debt-oriented Analysis.....</i>	<i>19</i>
1.3.2 <i>Sequence Classification to Predict Debt-related Activities.....</i>	<i>19</i>
1.3.3 <i>Combined Mining for Social Security Data.....</i>	<i>20</i>
1.4 THESIS CONTRIBUTION .....	20
1.5 THESIS ORGANIZATION .....	21
<b>2 REVIEW ON SOCIAL SECURITY/WELFARE DATA MINING RESEARCH .....</b>	<b>22</b>
2.1 INTRODUCTION .....	22
2.2 SOCIAL SECURITY SERVICES AND DATA.....	22

2.2.1	<i>Social Security Business</i>	22
2.2.2	<i>Social Security Data</i>	24
2.3	COMPREHENSIVE PICTURE	27
2.4	TECHNICAL PERSPECTIVE	28
2.5	RELATED WORK OF SOCIAL SECURITY DATA MINING	30
2.6	RETROSPECTION ON MINING SOCIAL SECURITY DATA	30
<b>3</b>	<b>SOCIAL SECURITY DATA MINING FRAMEWORK</b>	<b>33</b>
3.1	A BASIC FRAMEWORK	33
3.2	SOCIAL SECURITY DATA MINING GOALS	35
3.3	SOCIAL SECURITY DATA MINING CHALLENGES	37
3.3.1	<i>Social Security Data Processing</i>	37
3.3.2	<i>Social Security Pattern Analysis</i>	39
3.3.3	<i>Knowledge Delivery</i>	40
3.4	SOCIAL SECURITY DATA MINING TASKS	42
3.4.1	<i>Data Processing</i>	42
3.4.2	<i>Activity Analysis</i>	43
3.4.3	<i>Customer Risk Analysis</i>	44
3.4.4	<i>Earnings Analysis</i>	45
3.4.5	<i>Change Detection</i>	45
3.4.6	<i>Payment Accuracy Analysis</i>	46
3.4.7	<i>Fraud Detection</i>	48
<b>4</b>	<b>POSITIVE/NEGATIVE SOCIAL SECURITY SEQUENTIAL RULES ANALYSIS</b>	<b>49</b>

4.1	POSITIVE/NEGATIVE SEQUENTIAL RULES .....	49
4.1.1	<i>Introduction</i> .....	49
4.1.2	<i>Background and Related Work</i> .....	50
4.1.3	<i>Problem Statement</i> .....	52
4.2	EFFICIENT MINING OF EVENT-ORIENTED NEGATIVE SEQUENTIAL RULES .....	54
4.2.1	<i>Algorithm for Mining Event-Oriented Negative Sequential Rules</i> .....	54
4.2.2	<i>Experimental Evaluation</i> .....	56
4.2.3	<i>Conclusion</i> .....	59
4.3	MINING BOTH POSITIVE AND NEGATIVE IMPACT-ORIENTED SEQUENTIAL RULES FROM TRANSACTIONAL DATA60	
4.3.1	<i>Mining Impact-Oriented Sequential Rules</i> .....	60
4.3.2	<i>Experimental Results</i> .....	62
4.3.3	<i>Conclusions</i> .....	66
<b>5</b>	<b>PREDICT DEBT-RELATED SOCIAL SECURITY ACTIVITY SEQUENCES.....</b>	<b>67</b>
5.1	PROBLEM STATEMENT OF SEQUENCE CLASSIFICATION .....	67
5.2	SEQUENCE CLASSIFICATION USING BOTH POSITIVE AND NEGATIVE SEQUENTIAL PATTERNS .....	68
5.2.1	<i>Discriminative Sequential Patterns</i> .....	68
5.2.2	<i>Building Sequence Classifiers</i> .....	69
5.2.3	<i>Case Study</i> .....	70
5.3	DEBT DETECTION IN SOCIAL SECURITY BY ADAPTIVE SEQUENCE CLASSIFICATION.....	76
5.3.1	<i>Discriminative Frequent Patterns Boosting</i> .....	76
5.3.2	<i>Adaptive Sequence Classification Framework</i> .....	78
5.3.3	<i>Case Study</i> .....	80

5.4	CUSTOMER ACTIVITY SEQUENCE CLASSIFICATION FOR DEBT PREVENTION IN SOCIAL SECURITY .....	85
5.4.1	<i>Interestingness Measure</i> .....	85
5.4.2	<i>Sequence Classification</i> .....	86
5.4.3	<i>Case Study</i> .....	90
<b>6</b>	<b>MINING COMBINED SOCIAL SECURITY PATTERNS.....</b>	<b>94</b>
6.1	INTRODUCTION .....	94
6.2	RELATED WORK .....	94
6.3	THE PROBLEM .....	95
6.4	COMBINED PATTERN MINING .....	98
6.4.1	<i>Definitions of Combined Patterns</i> .....	98
6.4.2	<i>Interestingness Measures for Combined Patterns</i> .....	99
6.4.3	<i>Redundancy in Combined Patterns</i> .....	102
6.5	A CASE STUDY .....	103
<b>7</b>	<b>RARE CLASS ASSOCIATION RULE MINING WITH MULTIPLE IMBALANCED ATTRIBUTES.....</b>	<b>107</b>
7.1	BACKGROUND.....	107
7.1.1	<i>Class Association Rules</i> .....	108
7.1.2	<i>Data Imbalance in Association Rule Mining</i> .....	109
7.2	NOVEL ASSOCIATION RULE MINING PROCEDURE.....	109
7.2.1	<i>Interestingness Measures</i> .....	110
7.2.2	<i>Transformation</i> .....	112
7.3	TEST CASE .....	113
7.3.1	<i>Datasets Involved</i> .....	113

7.3.2	<i>Experimental Results</i> .....	114
7.4	CONCLUSIONS .....	117
<b>8</b>	<b>CONCLUSIONS AND FUTURE WORK</b> .....	<b>118</b>
8.1	CONCLUSIONS .....	118
8.2	FUTURE WORKS.....	118
8.2.1	<i>Data Mining Applications in Social Security</i> .....	118
8.2.2	<i>Sequential Rules Mining and Sequence Classification</i> .....	119
8.2.3	<i>Development of Further Models</i> .....	119
8.2.4	<i>A Straightforward Approach to Ongoing Research and Development</i> .....	120
	<b>LIST OF PUBLICATIONS</b> .....	<b>121</b>
	AWARDS .....	121
	BOOK CHAPTERS .....	121
	CONFERENCE PAPERS .....	121
	JOURNAL ARTICLES .....	122
	<b>REFERENCES</b> .....	<b>123</b>



## LIST OF FIGURES

FIGURE 2-1 SOCIAL SECURITY BUSINESS WORKFLOW .....	23
FIGURE 3-1 SSDM FRAMEWORK.....	33
FIGURE 3-2 SSDM GOALS .....	35
FIGURE 3-3 SSDM CHALLENGES.....	37
FIGURE 3-4 SSDM TASKS .....	42
FIGURE 4-1 ALGORITHM FOR DISCOVERING EVENT-ORIENTED NEGATIVE SEQUENTIAL RULES .....	56
FIGURE 4-2 SCALABILITY WITH MINIMUM SUPPORT.....	57
FIGURE 4-3 SCALABILITY WITH THE NUMBER OF SEQUENCES .....	58
FIGURE 4-4 SCALABILITY WITH THE NUMBER OF ITEMS PER SEQUENCE .....	58
FIGURE 4-5 SCALABILITY WITH THE AVERAGE LENGTH OF PATTERNS .....	59
FIGURE 4-6 PSEUDOCODE FOR DISCOVERING IMPACT-ORIENTED NEGATIVE SEQUENTIAL RULES.....	61
FIGURE 4-7 SCALABILITY WITH (A) SUPPORT; (B) THE NUMBER OF SEQUENCES; AND (C) THE LENGTH OF SEQUENCES.....	63
FIGURE 4-8 A GROWING SEQUENTIAL PATTERN “ADV ADV CCO” .....	66
FIGURE 5-1 ARCHITECTURE OF ADAPTIVE SEQUENCE CLASSIFICATION .....	79
FIGURE 5-2 ADAPTIVE CLASSIFICATION MODEL .....	80
FIGURE 5-3 EFFECTIVENESS OF DISCRIMINATIVE PATTERNS BOOSTING.....	82
FIGURE 5-4 ROC CURVES OF ADAPTIVE SEQUENCE CLASSIFICATION FRAMEWORK.....	84
FIGURE 5-5 HIERARCHICAL SEQUENCE CLASSIFICATION ALGORITHM.....	88
FIGURE 7-1 PROPOSED ALGORITHM .....	110
FIGURE 7-2 DISTRIBUTION OF THE IMBALANCED ATTRIBUTES.....	115

## LIST OF TABLES

TABLE 2-1 CENTRELINK BUSINESS DIMENSIONS 2008–2009 .....	24
TABLE 4-1 NOTATIONS .....	50
TABLE 4-2 SUPPORTS, CONFIDENCES AND LIFTS OF FOUR TYPES OF SEQUENTIAL RULES.....	54
TABLE 4-3 SELECTED POSITIVE AND NEGATIVE SEQUENTIAL RULES.....	65
TABLE 5-1 FEATURE-CLASS CONTINGENCY TABLE .....	69
TABLE 5-2 EXAMPLES OF ACTIVITY TRANSACTION DATA .....	71
TABLE 5-3 SELECTED POSITIVE AND NEGATIVE SEQUENTIAL RULES.....	72
TABLE 5-4 THE NUMBER OF PATTERNS IN PS10 AND PS05.....	73
TABLE 5-5 CLASSIFICATION RESULTS WITH PATTERN SET PS05-4K.....	75
TABLE 5-6 CLASSIFICATION RESULTS WITH PATTERN SET PS05-8K.....	75
TABLE 5-7 CLASSIFICATION RESULTS WITH PATTERN SET PS10-4K.....	75
TABLE 5-8 CLASSIFICATION RESULTS WITH PATTERN SET PS10-8K.....	76
TABLE 5-9 THE NUMBER OF PATTERNS IN THE FOUR PATTERN SETS.....	76
TABLE 5-10 CENTRELINK DATA SAMPLE .....	81
TABLE 5-11 DATA WINDOWS .....	83
TABLE 5-12 2 BY 2 FEATURE-CLASS CONTINGENCY TABLE.....	86
TABLE 5-13 SAMPLES OF CENTRELINK ACTIVITY DATA .....	91
TABLE 5-14 PERFORMANCE OF DIFFERENT ALGORITHMS .....	92
TABLE 5-15 COMPARISON OF THE PROPOSED ALGORITHM TO CONVENTIONAL ALGORITHM..	93
TABLE 6-1 TRANSACTIONAL DATA.....	96
TABLE 6-2 CUSTOMER DEMOGRAPHIC DATA .....	96
TABLE 6-3 TRADITIONAL ASSOCIATION RULES .....	96

TABLE 6-4 COMBINED ASSOCIATION RULES .....	96
TABLE 6-5 COMBINED RULE PAIRS .....	97
TABLE 6-6 TRADITIONAL ASSOCIATION RULES .....	104
TABLE 6-7 SELECTED COMBINED RULES .....	104
TABLE 6-8 SELECTED COMBINED RULE CLUSTERS.....	105
TABLE 7-1 SELECTED RESULTS WITH BALANCED ATTRIBUTES.....	114
TABLE 7-2 SELECTED RULES WITH IMBALANCED ATTRIBUTES CAPTION STYLE .....	116
TABLE 7-3 SELECTED RESULTS OF THE COMBINED ASSOCIATION RULES .....	116

## Abstract

Data mining in business applications has become an increasingly recognized and accepted area of enterprise data mining in recent years. In general, while the general principle and methodologies of data mining and machine learning are applicable for any business applications, it is often essential to develop specific theories, tools and systems for mining data in a particular domain such as social security and social welfare business. This necessity has led to the concept of *social security and social welfare data mining*, the focus of this thesis work.

Social security and social welfare business involves almost every citizen's life at different life periods. It provides fundamental and crucial government services and support to varied populations of specific need. A typical scenario in Australia is that it not only connects one third of our populations, but also associates with many relevant stakeholders, including banking business, taxation and Medicare. Such business engages complicated infrastructure, networks, mechanisms, policies, activities, and transactions. Data mining of such business is a brand new application area in the data mining community.

Mining such social welfare business and data is challenging. The challenges come from the unavailable benchmark and experience in the data mining for this particular domain, the complexities of social welfare business and data, the exploration of possible doable tasks, and the implementation of data mining techniques in relation to the business objectives.

In this thesis, which adopts a practice-based innovative attitude and focusses on the marriage of social welfare business with data mining, we believe we have realised our objective of providing a systematic and comprehensive overview of the social security and social welfare data mining. The main contributions consist of the following aspects:

- As the first work of its kind, to the best of our knowledge, we present an overall picture of social security and social welfare data mining, as a new domain driven data mining application.
- We explore the business nature of social security and social welfare, and the characteristics of social security data.

- We propose a concept map of social security data mining, catering for main complexities of social welfare business and data, as well as providing opportunities for exploring new research issues in the community.
- Several case studies are discussed, which demonstrate the technical development of social security data mining, and the innovative applications of existing data mining techniques.

The nature of social welfare is spreading widely across the world in both developed and developing countries. This thesis work therefore is timely and could be of important business and government value for better understanding our people, our policies, our objectives, and for better services of those people of genuine needs.